

P&S Modern SSDs

Introduction to MQSim

Rakesh Nadig

Prof. Onur Mutlu

ETH Zürich

Fall 2022

30th November 2022

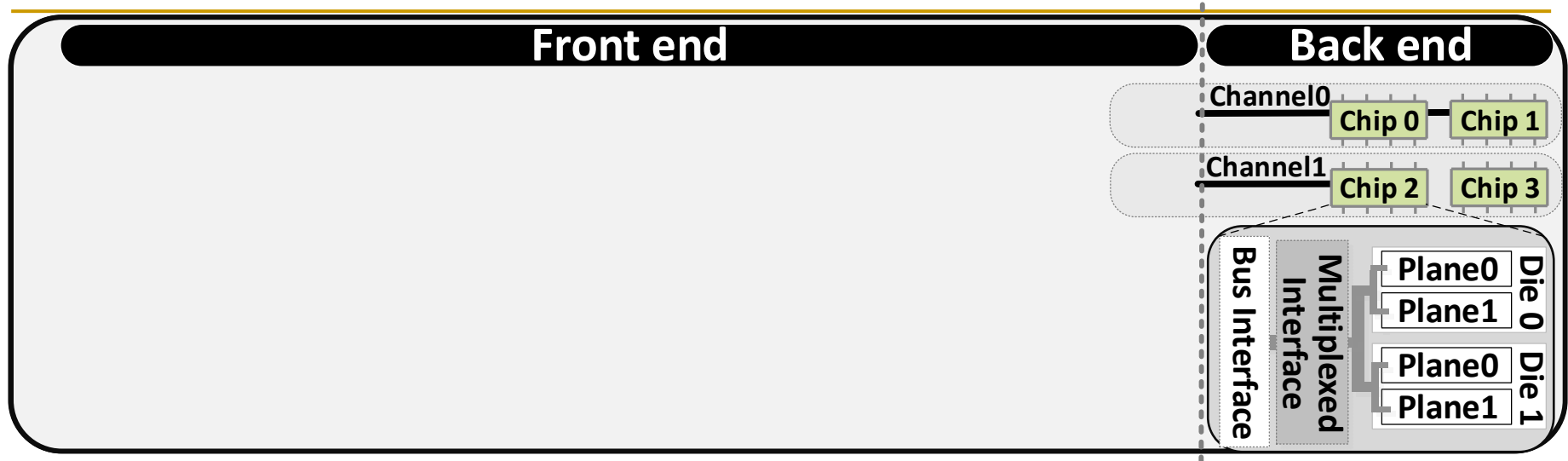
Outline

- Internal Components of a Modern SSD
- Introduction to MQSim
- Mechanism
- Code structure
- Configuring MQSim

Outline

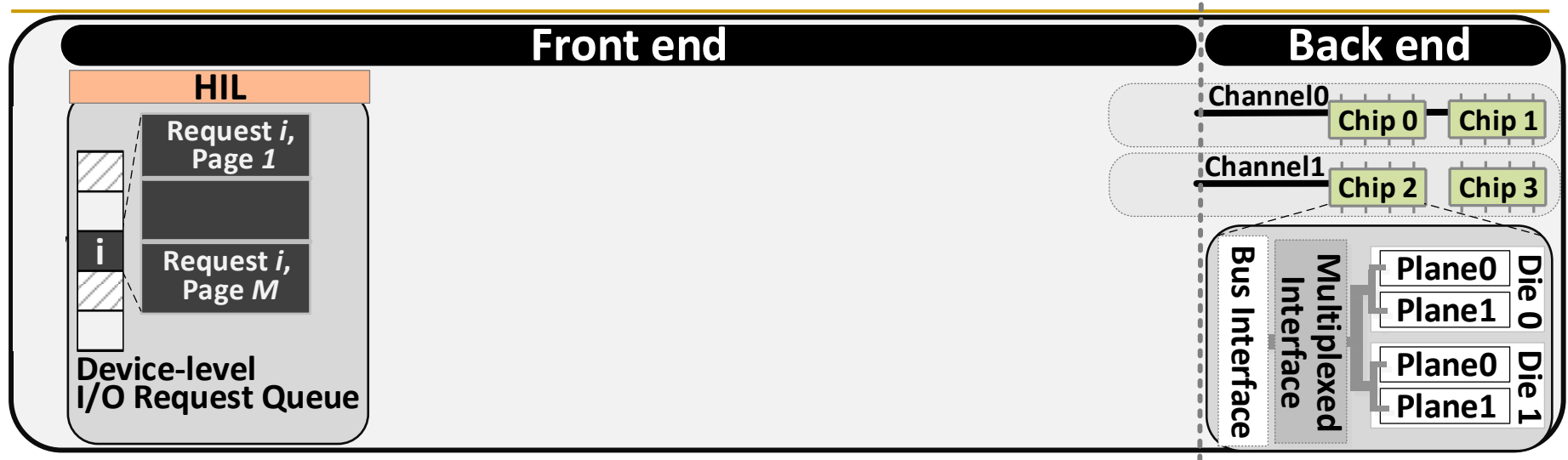
- Internal Components of a Modern SSD
- Introduction to MQSIM
- Mechanism
- Code structure
- Configuring MQSim

Internal Components of a Modern SSD



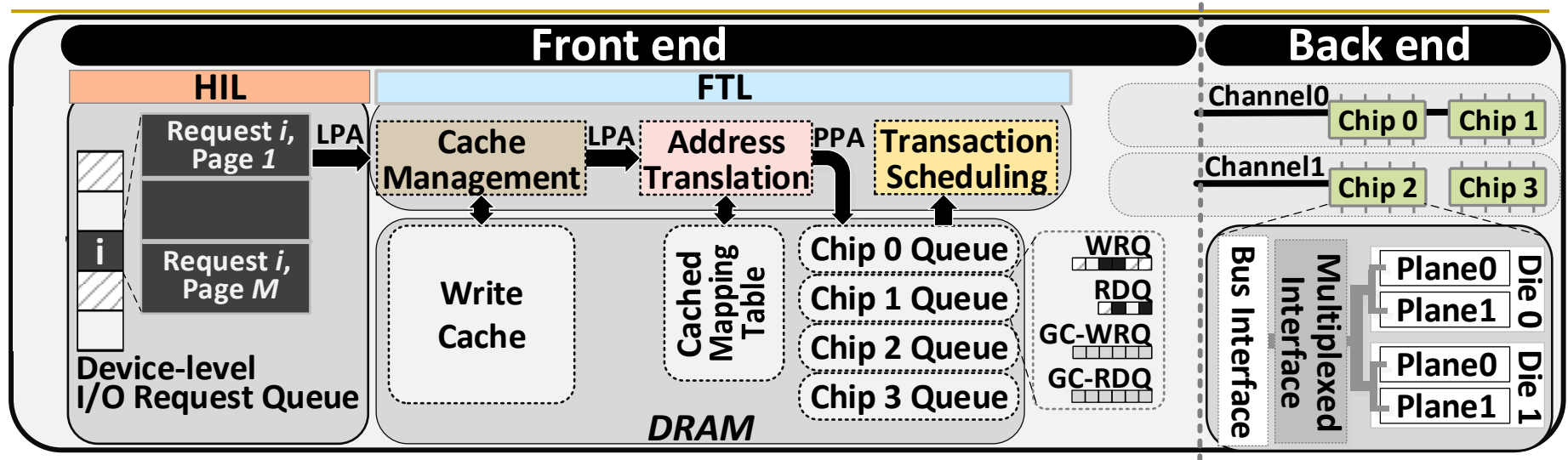
- Back End: data storage
 - Memory chips (e.g., NAND flash memory, PCM, MRAM, 3D XPoint)

Internal Components of a Modern SSD



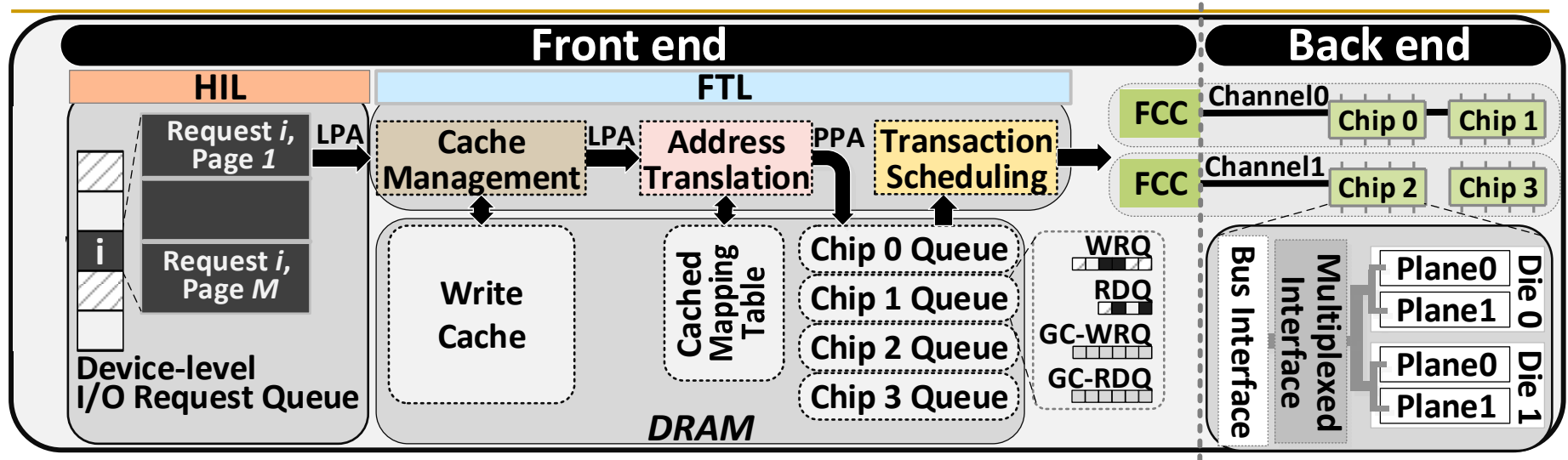
- Back End: data storage
 - Memory chips (e.g., NAND flash memory, PCM, MRAM, 3D XPoint)
- Front End: management and control units
 - **Host–Interface Logic (HIL):** protocol used to communicate with host (e.g., SATA, NVMe)

Internal Components of a Modern SSD



- Back End: data storage
 - ❑ Memory chips (e.g., NAND flash memory, PCM, MRAM, 3D XPoint)
- Front End: management and control units
 - ❑ **Host–Interface Logic (HIL):** protocol used to communicate with host (e.g., SATA, NVMe)
 - ❑ **Flash Translation Layer (FTL):** manages resources, processes I/O requests

Internal Components of a Modern SSD



- Back End: data storage

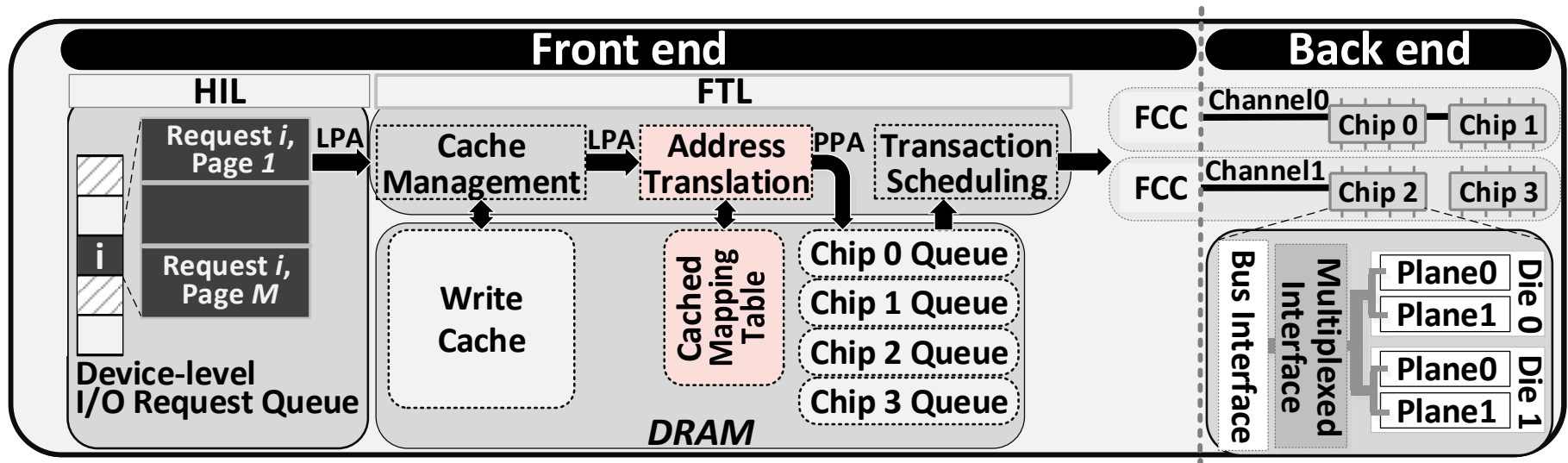
- Memory chips (e.g., NAND flash memory, PCM, MRAM, 3D XPoint)

- Front End: management and control units

- **Host-Interface Logic (HIL):** protocol used to communicate with host
- **Flash Translation Layer (FTL):** manages resources, processes I/O requests
- **Flash Channel Controllers (FCCs):** sends commands to memory chips, transfers data between the front end and back end

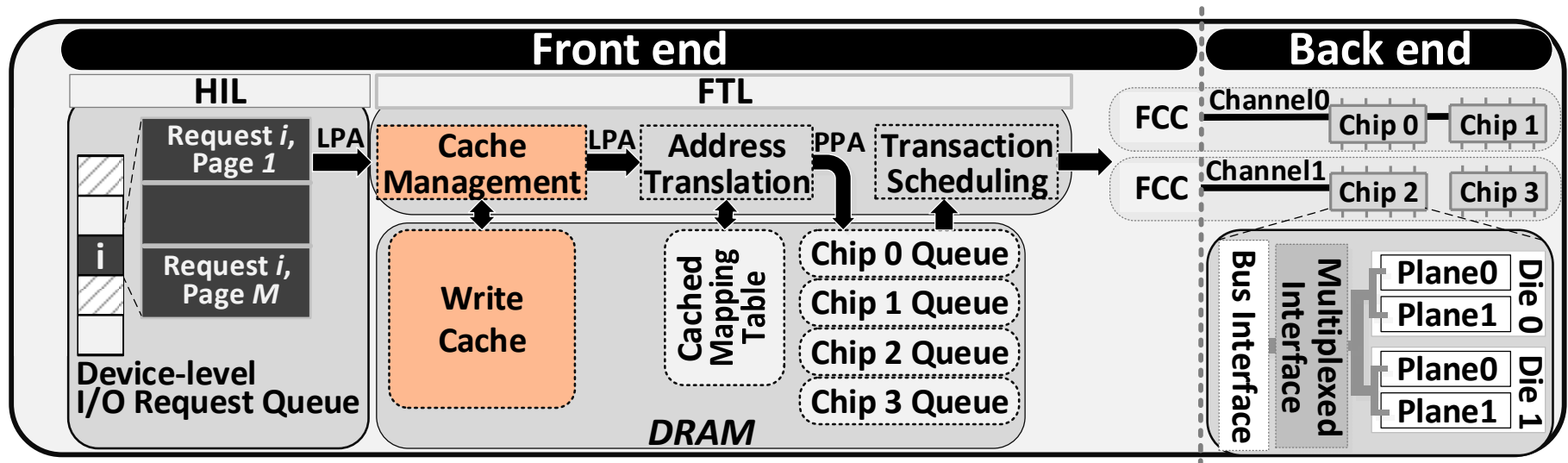
FTL: Managing the SSD's Resources

- Flash writes can be performed only on pages that are erased
 - Perform **out-of-place updates** (i.e., write data to a different, free page), mark **old page as invalid**
 - Update logical-to-physical mapping** (makes use of *cached mapping table*)
 - Some time later: **garbage collection** reclaims invalid physical pages *off the critical path of latency*



FTL: Managing the SSD's Resources

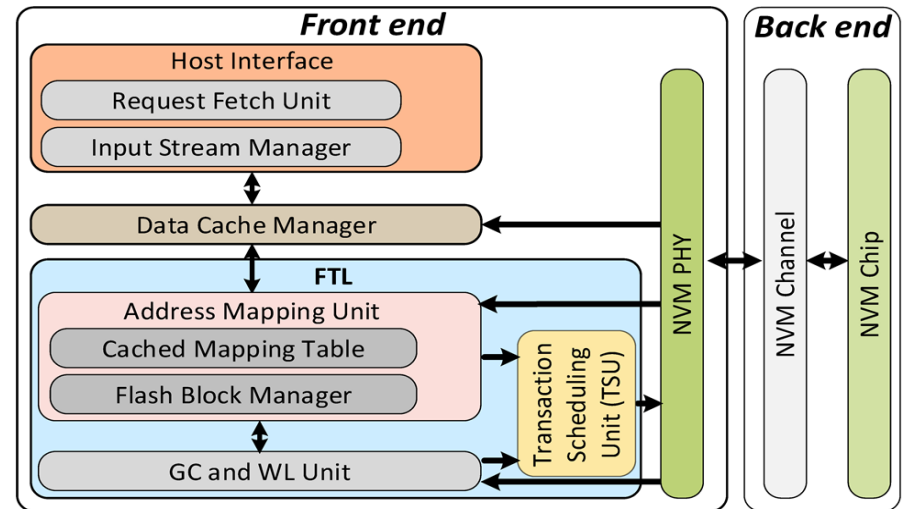
- Flash writes can take place only to pages that are erased
 - Perform **out-of-place updates** (i.e., write data to a different, free page), mark **old page as invalid**
 - Update logical-to-physical mapping** (makes use of **cached mapping table**)
 - Some time later: **garbage collection** reclaims invalid physical pages *off the critical path of latency*
- Write cache decreases resource contention, reduces latency



Outline

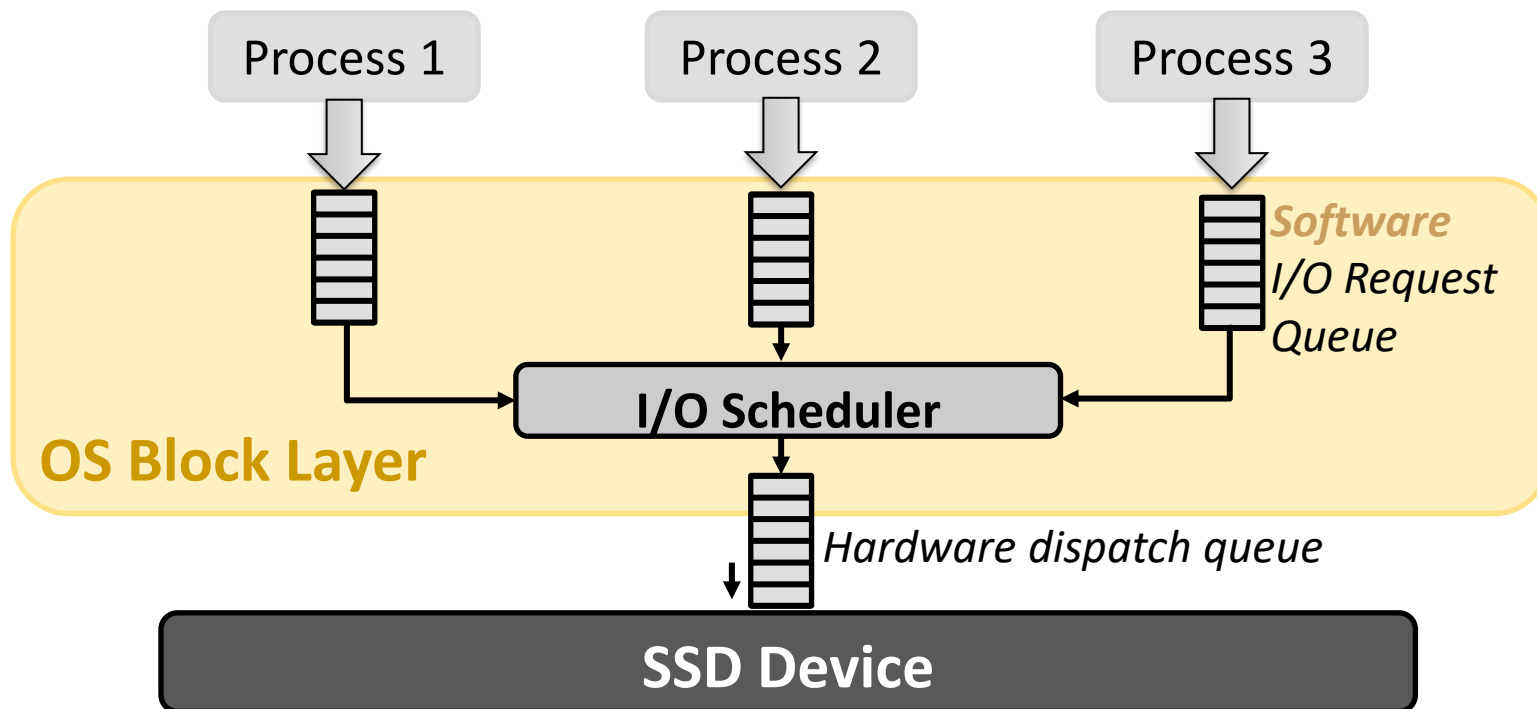
- Internal Components of a Modern SSD
- Introduction to MQSIM
- Mechanism
- Code structure
- Configuring MQSim

- Accurately models **conventional SATA-based SSDs** and **modern multi-queue SSDs**
 - ❑ Multi-queue protocols
 - ❑ Supports steady-state behavior with preconditioning
 - ❑ Models end-to-end I/O request latency
- Flexible design
 - ❑ Modular components
 - ❑ Ability to support **emerging non-volatile memory (NVM) technologies**
- Open-source release: <http://github.com/CMU-SAFARI/MQSim>
 - ❑ Written in C++
 - ❑ MIT License



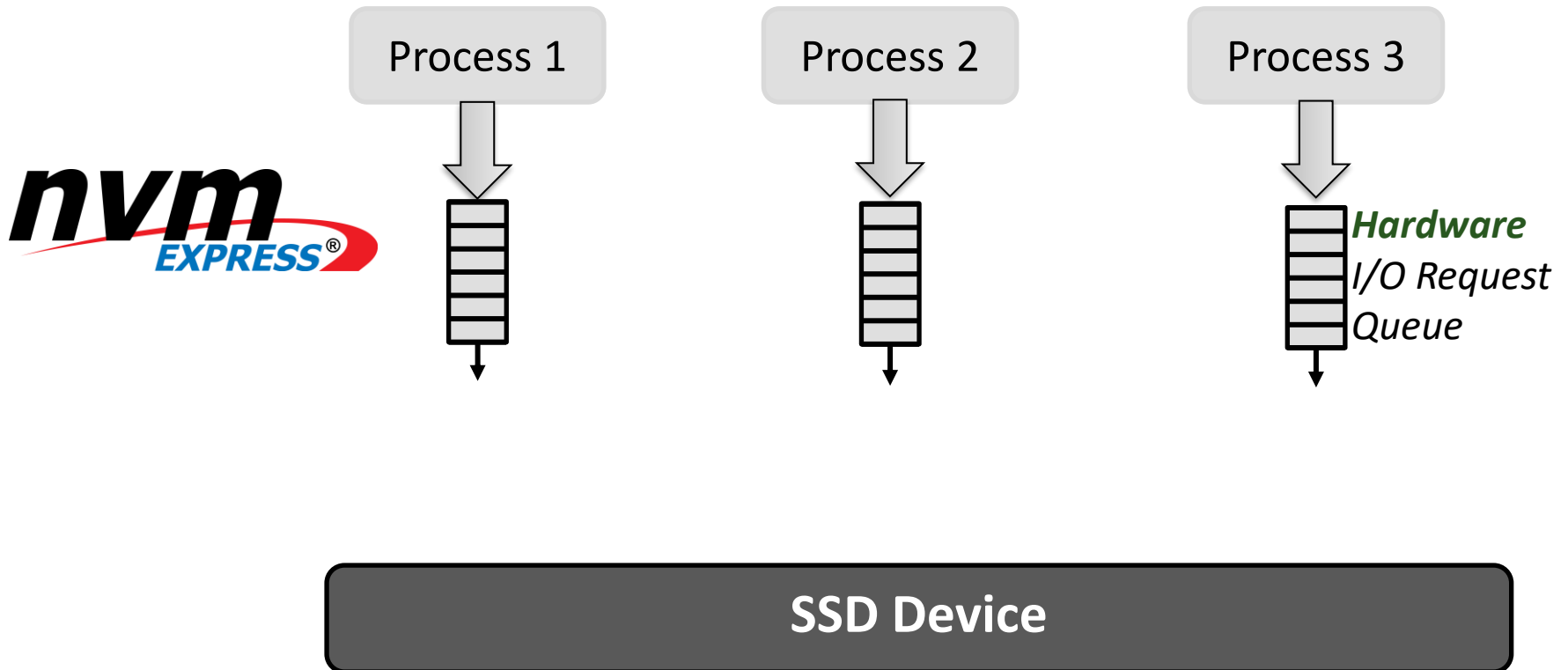
Support for Multi-Queue Protocols

- Conventional host interface (e.g., SATA)
 - ❑ Originally designed for magnetic hard disk drives: only **thousands of IOPS** per device
 - ❑ **OS handles scheduling, fairness** control for I/O requests



Support for Multi-Queue Protocols

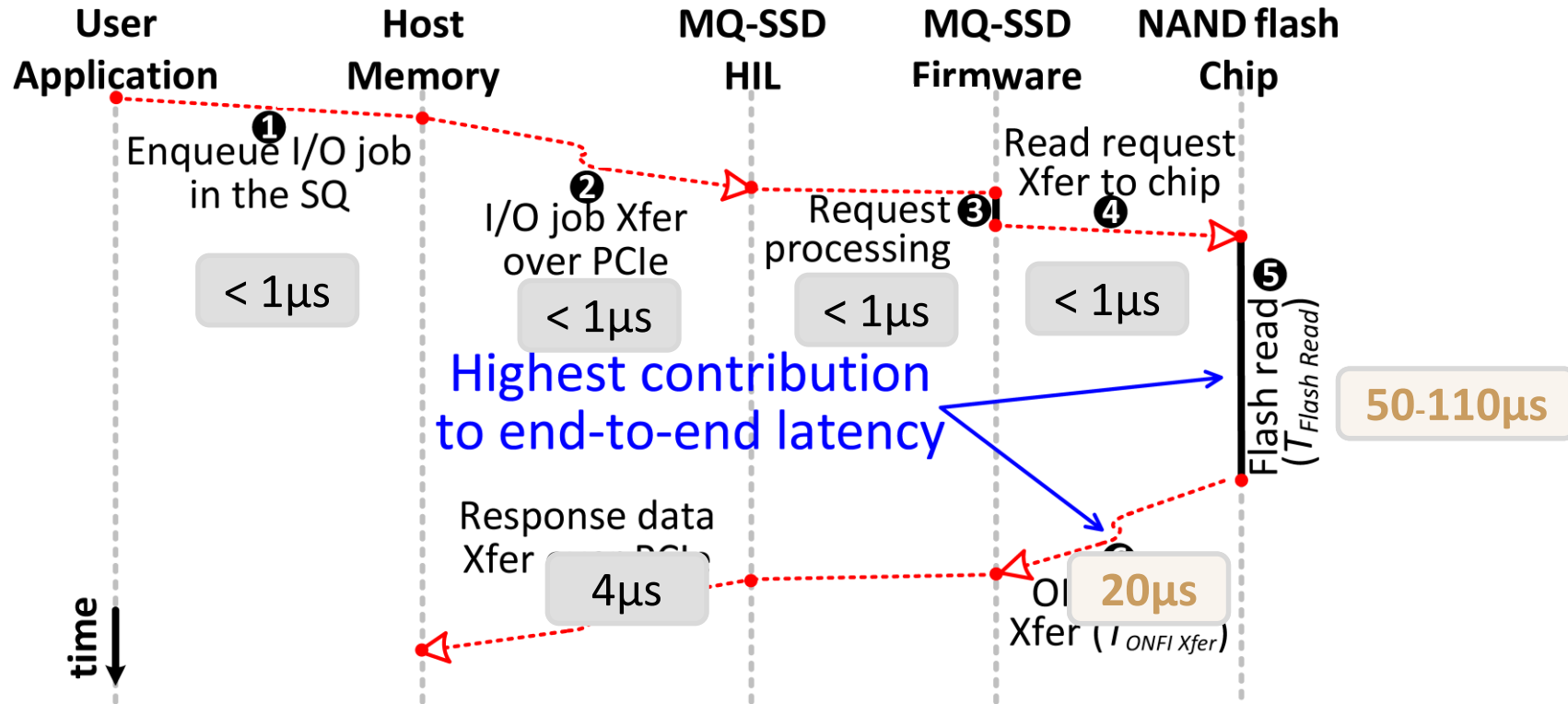
- Modern host interface (e.g., NVMe)
 - Takes advantage of SSD throughput: enables **millions of IOPS** per device
 - Avoids OS intervention: **SSD must perform scheduling, ensure fairness**



High-Performance Steady-State Model

- SSDs should be evaluated in steady state
 - *Fresh, out-of-the-box* (FOB) devices are **unlikely to perform garbage collection**
 - **Data cache is not warmed up** for an FOB device
- Many SSD simulators simulate only FOB devices or incorrectly model steady-state behaviour
- Difficult to reach steady state in most simulators
 - **Very slow**
 - **Widely-used traces aren't large enough for proper warm-up**

Complete Model of Request Latency



SSD Back End Model

- Three major latency components of the SSD back end
 - ❑ **Address and command transfer** to the memory chip
 - ❑ Flash memory **read/write execution** for different NAND technologies (1-bit, 2-bit and 3-bit per cell)
 - ❑ **Data transfer** to/from memory chips
- **Die- and plane-level parallelism** and advanced command execution
- Decouples the sizes of read and write operations
 - ❑ **Page-sized write** operations
 - ❑ **Sub-page read** operations

SSD Front End Model

- Host-Interface Model
 - ❑ **NVMe multi-queue (MQ)** and **SATA native command queue** models for a modern SSD
 - ❑ **Different priority classes** for **host side request queues** as per the NVMe standard specification
- Data Cache Manager
 - ❑ **DRAM-based cache** with Least Recently Used (LRU) replacement policy to cache recently accessed data
- FTL Components
 - ❑ Address translation unit
 - ❑ Garbage collection and wear-leveling unit
 - ❑ Transaction scheduling unit
 - ❑ Support for multi-flow request processing

More in the Paper

MQSim: A Framework for Enabling Realistic Studies of Modern Multi-Queue SSD Devices

Arash Tavakkol[†], Juan Gómez-Luna[†], Mohammad Sadrosadati[†], Saugata Ghose[‡], Onur Mutlu^{†‡}
[†]*ETH Zürich* [‡]*Carnegie Mellon University*

<https://www.usenix.org/system/files/conference/fast18/fast18-tavakkol.pdf>

Outline

- Internal Components of a Modern SSD
- Introduction to MQSIM
- **Mechanism**
- Code structure
- Configuring MQSim

Mechanism

- MQSim is a **discrete-event** based simulator
- Each operation performed by MQSim is an event (Command/Address transfer, Read data out etc.)
- A **red-black (RB)** tree maintains the events in the order of their completion times
- Each node of the RB tree has a **key** which indicates the **timestamp** to finish and, the **value** is the **callback** function
- For example, at time 5, a data transfer event starts, and the event will execute for another 5 time-units. The simulator will insert a node $\langle 5 + 5, \text{callback} \rangle$ into the tree
- At each simulation cycle, the simulator will find the node with the smallest time (value), set the timer to that timestamp, execute the callback function
- MQSim supports Real disk traces (e.g., MSR Cambridge Block I/O Traces, YCSB etc.) and Synthetic workloads (generated by the simulation engine)

Initialization

- **MQSim components** (FTL, SSD Device, TSU etc.) are represented as **objects** and maintained in a **map**
- The objects in the map are iterated through and corresponding functions are executed
- Each object maintains their own version of **setup_triggers**, **start_simulation** etc.

```
for(std::unordered_map<sim_object_id_type, Sim_Object*>::iterator obj = _ObjectList.begin();
    obj != _ObjectList.end();
    ++obj) {
    if (!obj->second->IsTriggersSetUp()) {
        obj->second->Setup_triggers();
    }
}

for (std::unordered_map<sim_object_id_type, Sim_Object*>::iterator obj = _ObjectList.begin();
    obj != _ObjectList.end();
    ++obj) {
    obj->second->Validate_simulation_config();
}

for (std::unordered_map<sim_object_id_type, Sim_Object*>::iterator obj = _ObjectList.begin();
    obj != _ObjectList.end();
    ++obj) {
    obj->second->Start_simulation();
}
```

Simulation Engine

■ Engine.cpp/start_simulation

```
while (true) {  
    if (_EventList->Count == 0 || stop) {  
        break;  
    }  
  
    EventTreeNode* minNode = _EventList->Get_min_node();  
    ev = minNode->FirstSimEvent;  
  
    _sim_time = ev->Fire_time;  
  
    while (ev != NULL) {  
        if(!ev->Ignore) {  
            ev->Target_sim_object->Execute_simulator_event(ev);  
        }  
        Sim_Event* consumed_event = ev;  
        ev = ev->Next_event;  
        delete consumed_event;  
    }  
    _EventList->Remove(minNode);  
}
```

The callback function is Execute_simulator_event

Mechanism of MQSim

- A simulator event is registered using the Register_sim_event function

```
Sim_Event* Engine::Register_sim_event(sim_time_type fireTime, Sim_Object* targetObject, void* parameters, int type)
{
    Sim_Event* ev = new Sim_Event(fireTime, targetObject, parameters, type);
    DEBUG("RegisterEvent " << fireTime << " " << targetObject)
    _EventList->Insert_sim_event(ev);
    return ev;
}
```

- Example: An event is registered for sending the Read command and address

```
if (chipBKE->OngoingDieCMDTransfers.size() == 0) {
    targetChip->StartCMDXfer();
    chipBKE->Status = ChipStatus::CMD_IN;
    chipBKE->Last_transfer_finish_time = Simulator->Time() + suspendTime + target_channel->ReadCommandTime[transaction_list.size()];
    Simulator->Register_sim_event(Simulator->Time() + suspendTime + target_channel->ReadCommandTime[transaction_list.size()], this,
        dieBKE, (int)NVDDR2_SimEventType::READ_CMD_ADDR_TRANSFERRED);
} else {
    dieBKE->DieInterleavedTime = suspendTime + target_channel->ReadCommandTime[transaction_list.size()];
    chipBKE->Last_transfer_finish_time += suspendTime + target_channel->ReadCommandTime[transaction_list.size()];
}
```

Outline

- Internal Components of a Modern SSD
- Introduction to MQSIM
- Mechanism
- **Code structure**
- Configuring MQSim

Code Structure

- **ssdconfig.xml** – configuration file that contains the preferred SSD configuration
- **workload.xml** – workload definition file
- traces – folder in which the trace files have to be placed
- src/sim – contains the files related to the simulation engine
- src/nvm_chip – code related to NVM chip. Contains files for die, plane, block and page
- src/ssd – contains files related to Flash Translation Layer and Transaction Scheduling Unit
- src/host – files related to trace and synthetic workloads, PCIe etc.

```
.
├── build
├── fast18
├── LICENSE
├── Makefile
├── MQSim
├── MQSim.exe
├── MQSim.pdb
├── MQSim.sln
├── MQSim.vcxproj
├── MQSim.vcxproj.filters
├── MQSim.vcxproj.user
├── README.md
├── src
├── ssdconfig.xml
├── traces
├── workload.xml
└── x64

├── src
│   ├── exec
│   ├── host
│   ├── main.cpp
│   ├── nvm_chip
│   ├── sim
│   ├── ssd
│   └── utils
├── ssdconfig.xml
├── traces
│   ├── tpcc-small.trace
│   └── wsrch-small.trace
└── workload.xml
```

Code Structure

- Code flow
 - Host -> Data Cache -> Address Translation Unit -> TSU -> Flash Controller -> NVM_Chip
- Some important functions:
 - Data Cache Manager: **process_new_user_requests()**
 - Address Mapping Unit: **translate_lpa_to_ppa_and_dispatch()**
 - Transaction Scheduling Unit: **Schedule()**
 - Flash Controller: **send_command_to_chip()**

Outline

- Internal Components of a Modern SSD
- Introduction to MQSIM
- Mechanism
- Code structure
- **Configuring MQSim**

SSD Configuration File (ssdconfig.xml)

- Host parameters
 - ❑ **PCIE_Lane_Bandwidth:** the PCIe bandwidth per lane in GB/s
 - ❑ **PCIE_Lane_Count:** the number of PCIe lanes
- Device Parameters
 - ❑ **HostInterface_Type:** the type of host interface. Range = {NVME, SATA}
 - ❑ **IO_Queue_Depth:** the length of the host-side I/O queue.
 - ❑ **Data_Cache_Capacity:** the size of the DRAM data cache in bytes
 - ❑ **Ideal_Mapping_Table:** if mapping is ideal, all the mapping entries are found in the DRAM and there is no need to read mapping entries from flash
 - ❑ **Transaction_Scheduling_Policy:** the transaction scheduling policy that is used in the SSD back end. Range = {OUT_OF_ORDER as defined in the Sprinkler (Jung et.al., HPCA 2014), PRIORITY_OUT_OF_ORDER which implements OUT_OF_ORDER and NVMe priorities}

SSD Configuration File (ssdconfig.xml)

■ Device Parameters

- ❑ **Flash_Channel_Count:** the number of flash channels in the SSD back end
- ❑ **Flash_Channel_Width:** the width of each flash channel in bytes
- ❑ **Chip_No_Per_Channel:** the number of flash chips attached to each channel in the SSD back end

■ NAND Parameters

- ❑ **Flash_Technology:** Range = {SLC, MLC, TLC}.
- ❑ Page read latency for LSB, CSB and MSB pages (in nanoseconds)
- ❑ Page program latency (in nanoseconds)
- ❑ **Block_Erase_Latency:** erase latency in nanoseconds
- ❑ **Block_PE_Cycles_Limit:** the PE limit of each flash block

SSD Configuration File (ssdconfig.xml)

■ NAND Parameters

- ❑ **Die_No_Per_Chip:** the number of dies in each flash chip
- ❑ **Plane_No_Per_Die:** the number of planes in each die
- ❑ **Block_No_Per_Plane:** the number of flash blocks in each plane
- ❑ **Page_No_Per_Block:** the number of physical pages in each flash block
- ❑ **Page_Capacity:** the size of each physical flash page in bytes
- ❑ **Page_Metadata_Capacity:** the size of the metadata area of each physical flash page in bytes

Trace-based workload configuration

```
<IO_Scenario>
  <IO_Flow_Parameter_Set_Trace_Based>
    <Priority_Class>HIGH</Priority_Class>
    <Device_Level_Data_Caching_Mode>WRITE_CACHE</Device_Level_Data_Caching_Mode>
    <Channel_IDs>0,1,2,3,4,5,6,7</Channel_IDs>
    <Chip_IDs>0,1,2,3</Chip_IDs>
    <Die_IDs>0,1</Die_IDs>
    <Plane_IDs>0,1</Plane_IDs>
    <Initial_Occupancy_Percentage>70</Initial_Occupancy_Percentage>
    <File_Path>traces/tpcc-small.trace</File_Path>
    <Percentage_To_Be_Executed>100</Percentage_To_Be_Executed>
    <Relay_Count>1</Relay_Count>
    <Time_Unit>NANOSECOND</Time_Unit>
  </IO_Flow_Parameter_Set_Trace_Based>
</IO_Scenario>
```

Synthetic workload configuration

```
<IO_Flow_Parameter_Set_Synthetic>
  <Priority_Class>HIGH</Priority_Class>
  <Device_Level_Data_Caching_Mode>WRITE_CACHE</Device_Level_Data_Caching_Mode>
  <Channel_IDs>0,1,2,3,4,5,6,7</Channel_IDs>
  <Chip_IDs>0,1,2,3</Chip_IDs>
  <Die_IDs>0,1</Die_IDs>
  <Plane_IDs>0,1</Plane_IDs>
  <Initial_Occupancy_Percentage>75</Initial_Occupancy_Percentage>
  <Working_Set_Percentage>50</Working_Set_Percentage>
  <Synthetic_Generator_Type>QUEUE_DEPTH</Synthetic_Generator_Type>
  <Read_Percentage>100</Read_Percentage>
  <Address_Distribution>RANDOM_UNIFORM</Address_Distribution>
  <Percentage_of_Hot_Region>0</Percentage_of_Hot_Region>
  <Generated_Aligned_Addresses>true</Generated_Aligned_Addresses>
  <Address_Alignment_Unit>16</Address_Alignment_Unit>
  <Request_Size_Distribution>FIXED</Request_Size_Distribution>
  <Average_Request_Size>8</Average_Request_Size>
  <Variance_Request_Size>0</Variance_Request_Size>
  <Seed>6533</Seed>
  <Average_No_of_Reqs_in_Queue>16</Average_No_of_Reqs_in_Queue>
  <Intensity>32768</Intensity>
  <Stop_Time>10000000000</Stop_Time>
  <Total_Requests_To_Generate>0</Total_Requests_To_Generate>
</IO_Flow_Parameter_Set_Synthetic>
```


MQSim Output File

```
<Host>
  <Host.IO_Flow>
    <Name>Host.IO_Flow.Trace.traces/tpcc-small.trace</Name>
    <Request_Count>6999</Request_Count>
    <Read_Request_Count>4381</Read_Request_Count>
    <Write_Request_Count>2618</Write_Request_Count>
    <IOPS>6999.000000</IOPS>
    <IOPS_Read>4381.000000</IOPS_Read>
    <IOPS_Write>2618.000000</IOPS_Write>
    <Bytes_Transferred>59718656.000000</Bytes_Transferred>
    <Bytes_Transferred_Read>36315136.000000</Bytes_Transferred_Read>
    <Bytes_Transferred_Write>23403520.000000</Bytes_Transferred_Write>
    <Bandwidth>59718656.000000</Bandwidth>
    <Bandwidth_Read>36315136.000000</Bandwidth_Read>
    <Bandwidth_Write>23403520.000000</Bandwidth_Write>
    <Device_Response_Time>3458</Device_Response_Time>
    <Min_Device_Response_Time>5</Min_Device_Response_Time>
    <Max_Device_Response_Time>18316</Max_Device_Response_Time>
    <End_to_End_Request_Delay>3458</End_to_End_Request_Delay>
    <Min_End_to_End_Request_Delay>5</Min_End_to_End_Request_Delay>
    <Max_End_to_End_Request_Delay>18316</Max_End_to_End_Request_Delay>
  </Host.IO_Flow>
</Host>
```

MQSim usage

■ Linux

- ❑ `$ make`
- ❑ `$./MQSim -i <SSD Configuration File> -w <Workload Definition File>`

■ Windows

- ❑ Open the MQSim.sln solution file in MS Visual Studio 2017 or later.
- ❑ Set the Solution Configuration to Release (it is set to Debug by default).
- ❑ Compile the solution.
- ❑ Run the generated executable file (e.g., MQSim.exe) either in command line mode or by clicking the MS Visual Studio run button. Please specify the paths to the files containing the 1) SSD configurations, and 2) workload definitions.
- ❑ `$ MQSim.exe -i <SSD Configuration File> -w <Workload Definition File>`

P&S Modern SSDs

Introduction to MQSim

Rakesh Nadig

Prof. Onur Mutlu

ETH Zürich

Fall 2022

30th November 2022